

1 We thank the reviewers for their thoughtful comments and helpful suggestions, and respond to comments below. We  
2 are pleased that the reviews are enthusiastic. R1: “*Interesting theoretical connections of maximum expected hitting*  
3 *cost to potential based reward shaping . . . This work may be useful in practice to craft informative shaping rewards*”  
4 and “*good significance. A tighter upper bound to diameter under the same set of assumptions is introduced while the*  
5 *connections to potential based reward shaping are interesting. Furthermore, this paper may be useful in practice to*  
6 *craft better shaping rewards.*”; R2: “*The paper is well-written and provides an interesting insight on the fact that the*  
7 *reward function in itself should ideally be taken into account in complexity measures.*”; and R3: “*the paper improves*  
8 *the regret bound on a very important algorithm by simple analysis*”

9 **R1: “bugs” in designed rewards (minor comment 1)** – In the particular context, we meant that from the perspective of the  
10 reward designer, the specified reward function might not be consistent with the desired behaviors (see L171). An  
11 example is the commonly cited OpenAI blog post on a faulty reward function in a game called CoastRunners (regrettably  
12 we cannot include any external links in the rebuttal). The quotation marks are to indicate the colloquial use of the term  
13 bug. We appreciate your suggestion and we shall add an explicit reference in any future versions.

14 **R1: discussion of contributions (minor comments 2-4)** – We agree with your assessment that the discussion section can be  
15 enhanced by emphasizing the *formal* results we established in this work. In contrast, in prior works, PBRS’s learning  
16 efficiency was only supported by numerical evidence. We are also excited to mention some prospects to extend this  
17 work in both theory and practice in any future versions.

18 **R2: applicability of MEHC beyond UCRL2** – We are also curious about the same question for non-optimism-based algorithms.  
19 We plan to pursue it in future works.

20 **R2:  $r_{max}$  in the definition of MEHC** – This is an interesting suggestion and we will consider it seriously. One benefit we  
21 currently see in keeping  $r_{max}$  in the definition is to remind readers (and users) that we assume some knowledge of  
22  $r_{max}$  in the bounded MDP setting (see footnote 2). Note also that a regret bound has a *unit* of “rewards” and even in the  
23 case of the unitless diameter, the resulting regret bound would include  $r_{max}$ , e.g.,  $\tilde{O}(r_{max}DS\sqrt{AT})$  from the original  
24 UCRL2 analysis.<sup>1</sup>

25 **R3: MEHC and diameter** – We agree that in *some* MDPs the gap between these two quantities can be small (or zero).  
26 However, we think this new complexity measure is worth studying as it provides a valuable tool to study the impact of  
27 rewards on learning efficiency as we have shown.

28 **R3: limited impact of PBRS on MEHC** – We refrain from making a harsh judgment on the merit of PBRS because in this  
29 work, we focus on the average reward setting and UCRL2 whose regret scales with MEHC. It is conceivable that for a  
30 different setting and a different RL algorithm that does not scale with MEHC (or scales poorly with a larger exponent),  
31 such as SARSA with epsilon-greedy exploration as used in [NHR99, footnote 4], PBRS might create a larger impact on  
32 learning efficiency as you suggested. We do agree that the discussion section may be enhanced by pointing out this  
33 future research direction and by adding qualifications lest the claims sound exaggerated.

34 **R3: PBRS vs other techniques** – As noted in the paper, PBRS is restrictive as the shaped rewards and the original rewards  
35 are  $\Pi$ -equivalent (see L156). It is reasonable to expect a pair of non- $\Pi$ -equivalent rewards—under some other weaker  
36 notion of equivalence—to have a greater difference in their learning efficiencies (under some algorithm). Furthermore  
37 we want to remark that [SBC05] studies a different RL setting with *salient events* and a comprehensive comparison of  
38 different means to incorporate expert knowledge is beyond the scope of this work.

39 **R3: a proof sketch** – We tried to make the detailed proofs (included as an appendix in the supplementary material) easy  
40 to follow, but we agree that a proof sketch will further enhance the paper by providing more intuition to readers. In  
41 particular, as you have complimented, Theorem 2 is a very interesting result and our proof, instead of constructing the  
42 best/worst potentials, relies only on the definitions of MEHC and that PBRS does not change rewards on a loop.

---

<sup>1</sup>In [JOA10],  $r_{max}$  is assumed to be 1 reward-unit making the expression *look* unitless.